



# Fraud Detection Using Transaction and Customer Data

Olatunji Ahmed

# The Critical Need for Early Fraud Detection

Financial fraud represents one of the most significant threats facing modern financial institutions and their customers. The sophistication of fraudulent schemes continues to evolve, making detection increasingly challenging. Traditional rule-based systems often fail to identify emerging fraud patterns, resulting in substantial financial losses and eroded customer trust.

Early detection of fraudulent transactions is not merely advantageous—it is critical for effective risk mitigation. The ability to identify suspicious activity before significant damage occurs can save institutions millions while protecting customer assets and maintaining regulatory compliance.



## Machine Learning Approach

Advanced algorithms detect complex fraud patterns



## Transaction-Level Data

Real-time analysis of individual transactions



## Customer Information

Account history and behavioral patterns

# Research Objectives

This research project aims to develop a comprehensive, data-driven approach to fraud detection that combines the power of machine learning with rich customer and transaction data. The following objectives guide our methodology and define success criteria for the model.

01

---

## Build a Supervised Fraud Detection Model

Develop a classification system trained on labeled historical data to distinguish fraudulent from legitimate transactions with high accuracy

03

---

## Predict Fraud Risk for Unseen Transactions

Deploy the model to evaluate new, previously unseen transactions in real-time, providing probability-based fraud risk scores

02

---

## Combine Customer and Transaction Features

Create a holistic view by integrating behavioral patterns, account characteristics, and transaction-level attributes for comprehensive risk assessment

04

---

## Demonstrate Real-World Fraud Analytics Workflow

Establish an end-to-end framework from data preparation through model deployment that mirrors operational fraud detection systems

# Dataset Composition and Structure

The dataset comprises **1,000 observations** representing a diverse range of customer transactions. This synthetic dataset was carefully constructed to mirror real-world fraud detection scenarios while maintaining data privacy and security standards.

Each observation contains multiple dimensions of information that collectively provide a comprehensive view of transaction risk. The integration of customer-level and transaction-level features enables the model to identify subtle patterns indicative of fraudulent behavior.



## Transaction Features

- **Transaction Amount:** Dollar value of the transaction
- **Anomaly Score:** Statistical measure of deviation from normal behavior
- **Transaction Category:** Type of transaction (e.g., retail, transfer, withdrawal)
- **Timestamp:** Date and time of transaction occurrence



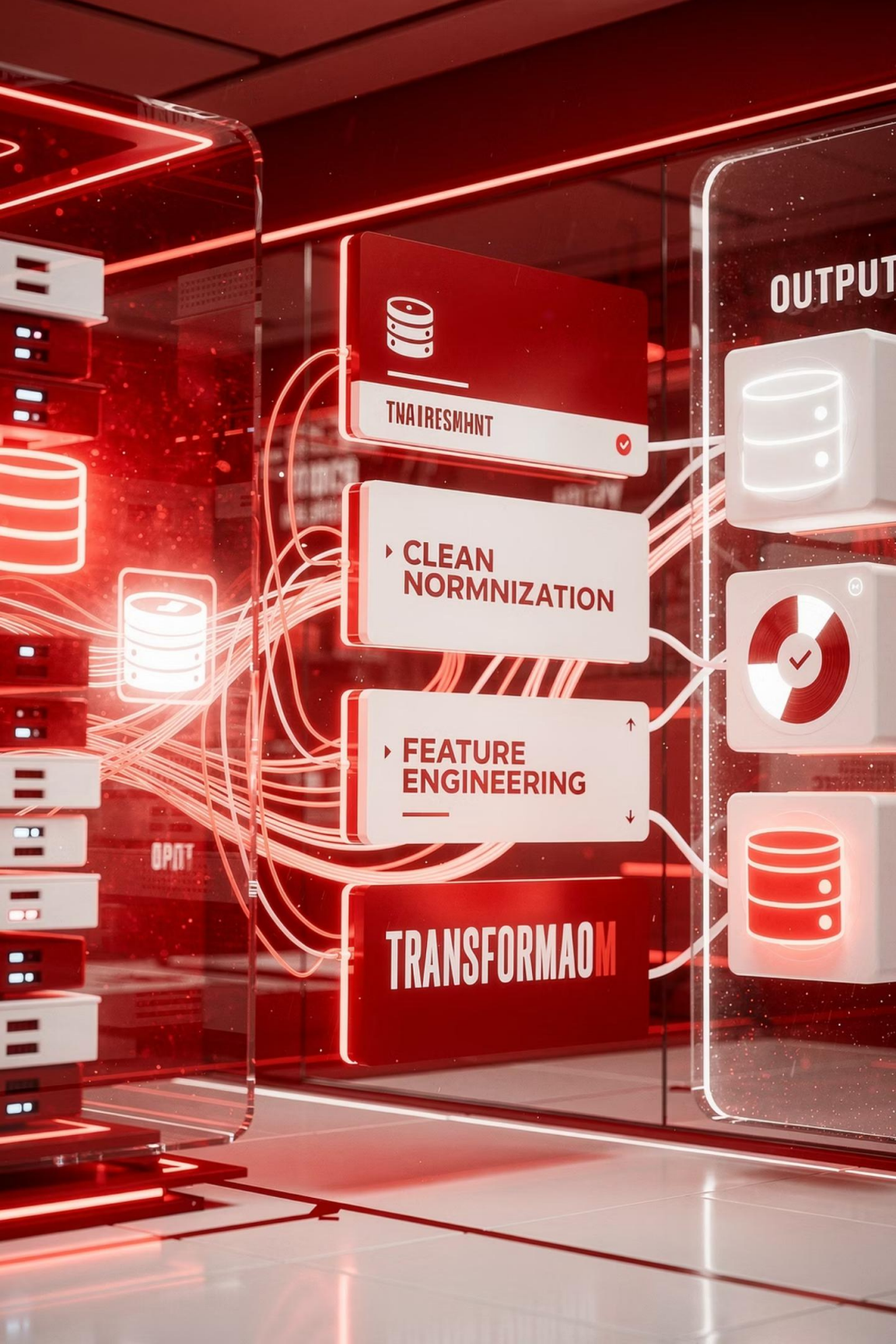
## Customer Features

- **Account Balance:** Current balance in customer account
- **Suspicious Flag:** Prior indicators of suspicious activity



## Target Variable

**FraudIndicator:** Binary classification where 0 indicates legitimate transactions and 1 indicates confirmed fraud cases



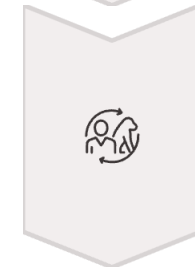
# Data Preparation and Preprocessing

Robust data preparation forms the foundation of any successful machine learning project. Our preprocessing pipeline ensures data quality, consistency, and proper formatting for model training while addressing common challenges in fraud detection datasets.



## Data Integration

Merged customer and transaction datasets on unique identifiers to create a unified analytical dataset



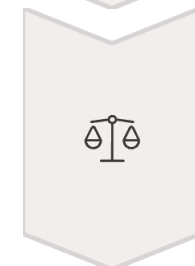
## Categorical Encoding

Applied One-Hot Encoding to transform categorical variables into numerical features suitable for machine learning algorithms



## Feature Alignment

Ensured consistent feature structure across training and inference datasets to prevent prediction errors



## Class Balancing

Implemented class weighting to address imbalanced fraud prevalence and improve minority class detection

# Feature Engineering Strategy

Feature engineering represents a critical phase where domain knowledge intersects with data science. The selected features capture two essential dimensions of fraud risk: **behavioral patterns** that indicate suspicious activity and **financial exposure** that quantifies potential loss.

Each feature was carefully chosen based on its predictive power and interpretability. The combination of categorical and continuous variables enables the model to detect both rule-based patterns and subtle statistical anomalies that might escape traditional detection systems.



Category

Type: Categorical  
Role: Transaction classification and pattern recognition

Transaction Amount

Type: Continuous  
Role: Primary indicator of financial exposure

Anomaly Score

Type: Continuous  
Role: Statistical measure of unusual behavior

Amount

Type: Continuous  
Role: Secondary transaction value metric

Account Balance

Type: Continuous  
Role: Context for transaction reasonableness

Suspicious Flag

Type: Binary  
Role: Historical behavior indicator



# Model Methodology and Implementation

## Supervised Learning Framework

Our approach employs a **supervised learning paradigm**, leveraging labeled historical data to train a classification model. This methodology allows the algorithm to learn patterns that distinguish fraudulent from legitimate transactions based on confirmed fraud cases.

## Logistic Regression Classifier

We selected **Logistic Regression** as our primary algorithm due to its interpretability, computational efficiency, and strong performance on binary classification tasks. The model outputs probability estimates rather than binary predictions, enabling risk-based decision thresholds.

### Scikit-learn Pipeline

Implemented using a production-ready pipeline that encapsulates preprocessing and model training in a single, reproducible workflow

### Automatic Preprocessing

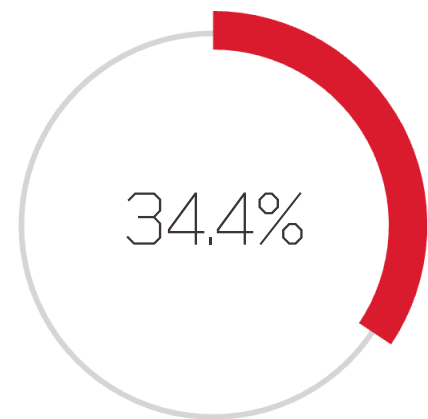
Embedded data transformations ensure consistent feature engineering across training and deployment environments

### Probability-Based Classification

Generates fraud probabilities for flexible threshold adjustment based on business risk tolerance

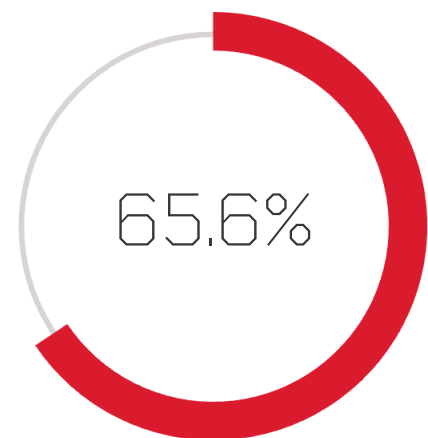
# Model Inference Results

To evaluate model performance, we analyzed an unseen transaction that was not part of the training dataset. This real-world test demonstrates the model's ability to generalize beyond historical examples and assess fraud risk for novel transactions.



Legitimate Transaction

Probability the transaction is not fraudulent



Fraudulent Transaction

Probability the transaction is fraud



## Classification Result: Potential Fraud

The model assigned a **fraud probability of 65.6%** to this transaction, which exceeds the standard decision threshold of 0.5 (50%). Based on this risk assessment, the transaction was flagged as potential fraud and would trigger additional verification procedures in a production environment.

This probability-based approach allows financial institutions to implement tiered response protocols: transactions with probabilities between 50-70% might receive automated verification, while those above 70% could trigger immediate manual review.



# Interpreting the Model's Prediction

It is essential to understand what the model's output represents and, equally important, what it does not represent. The fraud detection system **does not confirm fraud** with certainty—rather, it identifies **transactions with elevated fraud risk** that warrant further investigation.

This distinction mirrors the operational reality of fraud detection systems deployed by major financial institutions. These systems serve as the first line of defense, automatically screening millions of transactions and escalating high-risk cases to human analysts for final determination.

## Key Influencing Factors

Several features contributed significantly to the high fraud probability assigned to the test transaction:

- **High Anomaly Score:** The transaction exhibited statistical patterns that deviate substantially from the customer's historical behavior
- **Suspicious Behavior Indicators:** Previous flags in the customer's account history increased the baseline risk assessment
- **Transaction Characteristics:** The combination of transaction amount, category, and account balance created a risk profile consistent with known fraud patterns



“

"Effective fraud detection balances automation with human judgment. Machine learning identifies risk; human expertise confirms fraud."

”

# Study Limitations and Considerations

While this project successfully demonstrates a fraud detection workflow, several limitations must be acknowledged to properly contextualize the results and guide future research directions.



## Synthetic Dataset

The analysis relies on synthetic rather than real customer data. While designed to simulate realistic fraud scenarios, synthetic data may not capture the full complexity and nuance of actual fraudulent behavior patterns observed in production environments.



## Class Imbalance Effects

Despite implementing class weighting, the inherent imbalance between fraud and non-fraud cases (reflecting real-world distributions) may affect model sensitivity. The model might exhibit conservative behavior to maintain overall accuracy.



## Linear Relationship Assumption

Logistic regression assumes approximately linear relationships between features and the log-odds of fraud. Complex, non-linear interactions between variables may not be fully captured by this modeling approach.






## Limited Validation





Further validation with real-world datasets, temporal validation (testing on future data), and cross-institutional testing would be necessary before deployment in production fraud detection systems.

# Conclusion and Future Research Directions

## Key Achievements

-  **Successful Model Development**  
Built a functional fraud risk detection model capable of assessing transaction-level fraud probability with reasonable accuracy
-  **End-to-End ML Workflow**  
Demonstrated complete analytical pipeline from data preparation through model inference, establishing reproducible best practices
-  **Risk Identification Capability**  
Successfully identified high-risk transactions using previously unseen data, validating the model's generalization potential

## Future Research Opportunities

-  **Advanced Modeling Techniques**  
Test ensemble methods such as Random Forest and XGBoost to capture non-linear relationships and feature interactions
-  **Probability Calibration**  
Apply calibration techniques to ensure predicted probabilities accurately reflect true fraud likelihood
-  **Temporal Features**  
Incorporate transaction history and time-series patterns to enhance detection of evolving fraud schemes
-  **Production Deployment**  
Develop real-time monitoring system with streaming data processing for immediate fraud detection